

# Classifications and Grammars of Simple Arabic Lexical Invariants in Anticipation of an Automatic Processing of this Language “The Temporal Invariants”

Dhaou Ghoul<sup>1</sup>, André Jaccarini<sup>2</sup>, Amr Helmy Ibrahim<sup>1</sup>

<sup>1</sup> Sorbonne University,  
STIH Laboratory,  
France

<sup>2</sup> Aix-en-Provence University,  
MMSH,  
France

{dhaou.ghoul, amrhelmyibrahim, andre.jaccarini}@gmail.com

**Abstract.** Our study focuses on the classification and treatment of simple Arabic lexical invariants that express a temporal aspect. Our aim is to create a diagram of grammar (finite state machine) for each invariant. In this work, we limited our treatment to only 13 lexical invariants. Our hypothesis begins with the principle that the lexical invariants are located at the same structural level (formal) as the schemes in the language quotient (skeleton) of the Arabic language. They hide a great deal of information and involve syntactic expectations that makes it possible to predict the structure of the sentence. To do this: First, for each lexical invariant, we developed a sample corpus that contains the different contexts of the lexical invariant in question. Second, from this corpus, we identified the different linguistic criteria of the lexical invariant that allow us to correctly identify it. Finally, we codified this information in the form of linguistic rules in order to model it by a diagram of grammar (finite state machine).

**Keywords.** Corpus, classification, syntactic environment, regular expression, Arabic language, lexical invariants, identification, linguistic rules, regular grammar, diagrams of grammars, NLP.

## 1 Introduction

As part of Mogador<sup>1</sup> project [1] we are interested in tool words that we call lexical invariants. Indeed, morphosyntactic ambiguity is most often identified in the Arabic language. According to our research, we noticed that several arguments demonstrate

---

<sup>1</sup> Modélisation des grammaires arabes, des données et des outils de recherche. <https://halshs.archives-ouvertes.fr/halshs-00912009/document>.

the ambiguity of tokens: For example, various graphical shape due to the absence of vowels and treating the word out of context. By contrast, a mechanism for morphosyntactic disambiguation of Arabic is needed to resolve the ambiguity of unveiled words. The aim of our research is to treat the Arabic language with minimum although complex rules without referring to a lexicon. In their work Audebert *et al.* proposed grammars that define some syntactic operators in Arabic ('inna, 'anna, 'an...) [2]. Our work is based on this bibliography to improve the grammatical basis of the operators.

In general, lexical invariants are classified according to their meaning and their function in the sentence. Consequently, these invariants play an important role in the interpretation of the sentence and the coherence of a text.

Our linguistic study consists in treating the context of these lexical invariants based on a corpus that represent a sample of each lexical invariant.

This study is organized as follows: In section 2: we present the concept of "lexical invariant" by exposing the various levels of invariance. Then, we classify the simple lexical invariants according to several criteria. In section 3: we present our method of linguistic study, which includes modeling through diagrams of the grammar's some temporal lexical invariants. In section 4: we present two examples of simple lexical invariants. Conclusions and perspectives will be presented in section 5.

## 2 Definition and Classification of Simple Arabic Lexical Invariants

### 2.1 Definition of Lexical Invariant

As part of our research project, we try to project the Arabic language on its skeleton noted L/RAC. This skeleton is a quotient language whose elements are constituted of:

- Classes of equivalence of lexemes that are the schemes (patterns, lexical paradigm, wazn).
- Singletons (classes with only one element) which we can not associate a schemes (wazn). So, these singletons represent the invariants of the Arabic language that retain their same forms in the quotient language. This is why we have called "lexical invariants".

This language is an abstract object *-formally-* constructed from the free monoid constituted by the set of concatenations of the Arabic graphemes  $A^*$  (with,  $A$ : Arabic alphabet and  $*$  designating the Kleene star) which the language  $L_{AR}$  consisting of the set of licit. Arabic graphic forms (graphic words separated by two whites) is a strict subset. What we mean by LEX is a particular subset of  $L_{AR} \subset A^*$ . Thus,  $LEX \subset L_{AR} \subset A^*$ . As well as, LEX is obtained by systematic segmentation of the elements of  $L_{AR}$ . [3]. Indeed, our appellation of "lexical invariant" is inspired by the definition of André Jaccarini [4].

**Table 1.** Temporal lexical invariants distribution.

	#sentences	#words	#lexical invariant	#temporal invariant	%temporal invariant
Text1	142	3460	1174	173	14.3
Text2	63	1311	507	80	15.7
Text3	123	2366	752	109	14.49
Text4	11	344	108	10	9.25

Let SC be the morphism associated with the operation of projection on the scheme.

We call the lexical invariant every element  $x$  belongs to LEX that is invariant with respect to the morphism SC.

The set of lexical invariants is denoted by:

$$INL = LEX \text{ inter } \{x ; SC(x) = \{x\}\} \{x ; x \in LEX \text{ et } SC(x) = \{x\}\}. \quad (1)$$

For example, SC (فإنهم / *fa'inahum*) = ف SC (إن) هم = إن. ف. With إن  $\in$  LEX, So إن is a lexical invariant.

These lexical invariants play a crucial role because they are the only elements that are both in the terminal vocabulary of L and L/RAC. They are on the same level as the schemes in L/RAC (quotient language) and they escape from the rules of morphological derivation. With, L is the language which consists of the set of grammatical sentences - not necessarily endowed with meaning-.

In Arabic there are two types of lexical invariants: simple lexical invariants like "*hattā*" and complex lexical invariants like "*hīnamā*". The latter are on the form " $x \ m\bar{a}$ " (with  $x$  representing a simple lexical invariant). In this study we treat only the simple lexical invariants.

## 2.2 Classification of Simple Temporal Arabic Lexical Invariants

Temporal lexical invariant classification remains an indispensable step in the automatic processing of the Arabic language. Based on literature on the one hand and on our study of these invariants on the other hand, we have tried to classify these 13 simple temporal invariants. Moreover, we calculated the percentage of temporal invariant's occurrences with respect to the total number of lexical invariants in four texts selected via our working corpus, which we created in previous work [5, 6]. Our calculation is made with the application "*Kawākib*" [7] followed by a manual check.

The following table shows the temporal lexical invariants distribution with respect to all invariants in a given text.

The aim of this classification is to show the linguistic criteria of these types of words in Arabic and their influence in the construction of a text. The list of temporal invariants that we treat in this work is as follows: *l.m.ā*; *'aṭnā'a*; *'idā*; *'id*; *hattā*; *hīn/ hīna*; *lam*; *lan*; *munḍu*; *qad*; *tumma*; *ba'da*; *sawfa*.

Our classification of simple temporal lexical invariants is based on five linguistic criteria: agglutination, based on the word that follows them, syntactic role, ambiguity and based on their *rection* in the sentence.

**Table 2.** Lexical invariants classification based on their agglutination.

Class	Classification criteria	Example
1	Agglutination only to a coordinating conjunction.	<i>l.m.ā, 'iq, ḥattā</i>
2	Agglutination to a coordinating conjunction and / or an interrogative conjunction.	<i>'iqā, lam, lan</i>
3	Agglutination to a coordinating conjunction and / or a corroborating conjunction.	<i>qad</i>
4	Agglutination to a coordinating conjunctions and / or preposition and a personal pronoun attached.	<i>'aṭnā'a</i>
5	Agglutination to a coordinating conjunctions and / or preposition and a personal pronoun attached and to a relative pronoun.	<i>hīn</i>
6	Agglutination to a coordinating conjunction and / or a personal pronoun attached and a relative pronoun.	<i>ba'da</i>
7	Agglutination to a coordinating conjunction and / or an interrogative conjunction and a corroborating conjunction.	<i>Sawfa</i>

**Table 3.** Lexical invariants classification based on the word that follows them.

Class	Classification criteria	Example
1	Pro-nominal invariants	<i>'aṭnā'a, ba'da</i>
2	Pre-verbal invariants	<i>qad, lam, lan, sawfa</i>
3	Mixed invariants	<i>l.m.ā, ḥattā, 'iq, 'iqā, tumma, munḍu, hīna</i>

**Lexical invariants classification based on their agglutination** with the prefixes or suffixes, we classified these 13 simple lexical invariants into 7 classes as the following table 2 shows.

**Lexical invariants classification based on the word that follows them.** According to our study of lexical invariants, it seemed useful to classify an invariant based on the word that follows them. Therefore, based on this criterion we have classified these 13 simple lexical invariants into three classes or families: pro-nominal invariants<sup>2</sup>, pre-verbal invariants<sup>3</sup>, mixed invariants<sup>4</sup> as the following table 3 shows.

**Lexical invariants classification based on their syntactic roles.** In Arabic language, we can find lexical invariants whose graphic form is unique but can have several syntactic functions. Among these lexical invariants is “*ḥattā*” whose unique graphic form may contain a subordinating conjunction, coordinating conjunction, adverb or preposition. Based on their syntactic roles, these lexical invariants are divided into 9 classes as the following table 4 shows.

**Lexical invariant classification based on their “rection” in the sentence.** Some lexical invariants in Arabic change the vocalization of the word that follows them. This

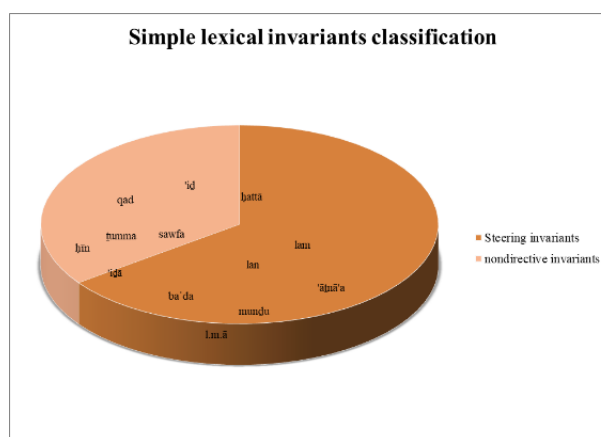
<sup>2</sup> The lexical invariants that precede the noun or the elements equivalent to the noun.

<sup>3</sup> The lexical invariants that precede the verb or the elements that can receive the verb.

<sup>4</sup> The lexical invariants that precede either a verb or a name or another invariant.

**Table 4.** Lexical invariants classification based on their syntactic roles.

Class	Syntactic class	Example
1	Subordinating conjunction	<i>lammā, lamā, limā, 'id, ḥattā, 'idā, 'aṭnā'a, ḥīna, munḍu</i>
2	Coordinating conjunction	<i>ḥattā, tumma</i>
3	Negation particle	<i>lan, lam, lammā</i>
4	Interrogative pronoun	<i>limā</i>
5	Preposition	<i>ḥīna, 'aṭnā'a, ba'd, ḥattā, munḍu</i>
6	Temporal adverb	<i>l.m.ā, 'id, 'idā, 'aṭnā'a, ḥīna, munḍu, ba'd, sawfa, ḥattā, lan, lam, tumma, qad</i>
7	Conditional conjunction	<i>'id, 'idā</i>
8	Surprise particle	<i>'id, 'idā</i>
9	Future particle	<i>sawfa</i>

**Fig. 1.** Lexical invariants classification based on their “reaction” in the sentence.

linguistic criterion is useful in the identification of these lexical invariants. It is therefore important to classify these invariants according to their power of reaction. Indeed, and based on this criterion, we classified these lexical invariants into two classes as the following figure (Fig.1) shows:

**Lexical invariant classification based on their ambiguity.** The vowels absence in the Arabic text is one of the major problems of automatic processing of this language. This problem leads to several cases of grammatical (syntactic) and graphical or even semantic ambiguity. Moreover, the phenomenon of agglutination in Arabic can be a major source of ambiguity. Therefore, it seemed useful to classify our study invariants according to their ambiguity (graphic, syntactic and semantic) as the following table 5 shows.

It is worth noting that the different abbreviations are illustrated in the appendix. In general, the lexical invariants classification depends on their linguistic properties. Therefore, this classification is always debatable and not exhaustive.

**Table 5.** Lexical invariants classification based on their ambiguity.

Lexical invariant	Graphic ambiguity	Syntactic ambiguity	Semantic ambiguity
<i>l.m.ā</i>	<i>lammā, limā, lamā</i>	CS, PNE, ADVT	Correspondence, temporal, negation
'iḍ	'iḍ, āḍ	CS, COND, ADVT, PS	Condition, surprise
'iḍā	'iḍā, āḍā, āḍan, 'iḍan	CS, PS, CCOND, ADVT	Condition surprise
lan	Unique	ADVT, PNE	Negation
lam	Unique	ADVT, PNE	Negation
ḥattā	Unique	CS, CC, PREP, ADVT	Coordination, intensity, justification, direction...
tumma	<i>tumma, tamma</i>	CC, ADVT, PD, Verb	Ranking (order)
mundu	Unique	CS, ADVT, PREP	Time report, temporal interrogation
qad	Unique	ADVT	Certainty, uncertainty
ḥīna	ḥīna, ḥīn	ADVT, CS, PREP	Temporal correspondence, a moment
'atnā'a	'atnā'a, atnā'	CS, ADVT, PREP	Temporal synchronization
	ba'da,	PREP	Temporal, Continuity
b'd	ba'du,	ADVT	Anteriority
	bu'd,	Noun	Distance
	ba'ida, ba'uda, ba'ada	Verb	Move away
sawfa	Unique	ADVT, PF	Promise

### 3 Methodology of Lexical Invariants Modeling by Grammars Diagram

We will start from the idea that a grammar is a point of view among others on the language and that these points of view are related to the tasks assigned to these grammars. The grammars of the invariants are made up on them and their linguistic analysis. Information about them is largely known by Arab and Western grammarians. However, our contribution is summarized in the angle of analysis that we adopt and the examples that we select.

So, the grammars presented in our work are intended for the automated morpho-syntactic analysis of the Arabic language. The idea here is to transform the linguistic rules of each invariant that we have defined previously into a non-fixed-form formal grammar. These grammars describe the syntactic expectations of invariants in Arabic [8].

#### 3.1 Different Levels of Our Method

Our method of Arabic lexical invariant modeling is located as part of the theory of "abstract machines". It is based on two major phases (two levels of analysis).

**The linguistic level.** Our methodology of the study linguistic of each lexical invariant consists to study its environment based on a sample of representative sentences. The different steps of our method are:

- Identification of lexical invariants in the corpus.
- Treatment of the environment (left and right context) of the lexical invariant.
- Develop hypotheses on the structure of the sentence based on the principle of the "empty sentence".
- Identify the different characteristics of the lexical invariants (grammatical role, ambiguity, suppressibility or not, graphical form, distributional analysis [9, 10], local and global scope...).
- Study of the syntactic expectations of lexical invariants.
- Codify the different information into non-fixed formal grammars (rules of production or linguistic rewriting), to bring out the relations in order to make them usable by the machine.

Note that the semantic part of each invariant will be studied in parallel and can also be related to the definition of discriminating criteria.

**The graphic level.** This level consists in modeling the grammars in the form of finite state automata (grammars schemes) in order to make it comprehensible by a machine. The graphs represent successions of lexical invariant's syntactic expectations. This means that one passes from a state represented by a node to another state by means of a transition which is symbolized by an arc. These arcs are labeled with the terminal symbols of the grammar.

So why use the theory of automata to represent or treat the Arabic language?

Indeed, the structural (strongly grammatical) characteristics of Arabic are also algorithmic characteristics. These characteristics are easily translatable as part of the theory of automata. Thus, the theory of automata offers the possibility of a good study of Arabic in the universe of knowledge.

### 3.2 Limits of Our Method

In a general way, our method, purely algorithmic, differs from those of contemporary research. It uses minimal resources and is independent of lexicons. The origin of this approach is due to the representation of the Arabic language in the form of a quotient or skeleton language. The latter can be obtained by reducing all the roots into a single control root "fa'ala" to represent all Arabic schemes.

Indeed, at this stage we have seen that in the language quotient exist particular words that we note by "lexical invariants". Our method consists in a surface analysis of these invariants in order to construct non-fixed formal grammars that represent the environment of each invariant. These grammars can be linked and combined. However, our method presents some limits:

- "Ad hoc" construction designed for a given task: The automata may contain unnecessary parts.

- Difficulty in marking the meaning of the lexical invariant in the automaton.
- Grammar diagrams may in some cases produce grammatical correct, but semantically incorrect sentences (invalid sentences because of the presence of a succession of epsilons).

## 4 Examples of Simple Temporal Lexical Invariants

In this section we will present two examples of simple lexical invariants among the 13 invariants that we have analyzed in this work. First of all, we will extract the maximum of information in the form of linguistic rules. Then, we will model these rules in regular grammars acceptable by finite state automata.

### 4.1 Lexical Invariant “*hīn /hīna*”

Outside of its context and in the absence of vowels, the word “*hīn*” is ambiguous grammatically. Indeed, this word may be a temporal subordinating conjunction or a noun. The disambiguation of this word depends on its position in the sentence as well as on the linguistic study of its environment. To assign the correct grammatical function to “*hīn*”, we must treat its context based on a sample corpus. This sample corpus represents an extract of 60 sentences from our working corpus [5, 6].

Independently of its grammatical function, the removal of “*hīn*” causes the disruption of the structure as well as the meaning of the sentence.

**Temporal subordinating conjunction.** In this case, the lexical invariant “*hīn*” expresses a relation of time between two events and noted “*hīna*”. This is a link between two parallel events as shown in the following example (1). In addition, it is very fertile lexical invariant. That is, it can agglutinate at  $n + 1$  with suffixes like “*mā*,” “*idan*,” “*dāk*” to establish a new lexical invariant.

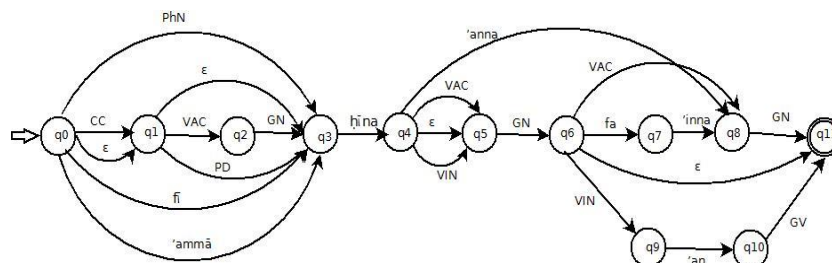
« أصابته الدهشة حين زار مدينة القاهرة » *'asābathu aldahṣat hīna zāra madīnat alqāhirat.*  
(He was astonished when he visited the city of Cairo). (1)

« أصابته الدهشة عندما زار مدينة القاهرة » *'asābathu aldahṣat 'indamā zāra madīnat alqāhirat.*  
» (He was astonished when he visited the city of Cairo). (1')

From the examples (1) and (1') above, we note that we can replace “*hīna*” by “*indamā*” without any change in the sentence either in the sense or in the structure. However, this permutation is only valid in the case where “*hīna*” has the same syntactic expectations (a verb phrase) as “*indamā*” (sentence 1'). This permutation is strictly forbidden in the case where “*hīna*” is followed by a nominal group as shown in the example (2):

« في حين أن الكثير من الشباب ملتزمون بجلسات المقاهي » *fī hīna 'anna alkaṭīr min alšabāb multazimūn bigalasāt almaqāhī.* (However, many young people were keen on going to coffees). (2)

If we look at the example (2), the sequence “*fī hīna*” is equivalent to “*baynamā*” (2') which marks a temporal opposition between two non-identical events. So, the lexical invariant “*hīna*” does not keep the same meaning when it appears in a sequence of invariants like “*fī hīna*”. It no longer expresses temporality, but it expresses an opposition.



**Fig. 2.** Grammar diagram of “*hīna*” temporal subordinating conjunction.

بينما أن الكثير من الشباب ملتزمون بجلسات المقاهي « *baynamā 'anna alkaṭīr min alšabāb multazimūn bigalasāt almaqāhī.* » (By contrast, many young people regularly attended cafés). (2')

As for syntactic expectations, the lexical invariant “*hīna*” in this case, waits in  $n + 1$  either a verb phrase (the verb can be either in the present tense or in the past tense), or a nominal group. However, it can be preceded in  $n - 1$  by a coordinating conjunction, the preposition “*fī*”, a demonstrative pronoun, corroborating conjunction “*ammā*” or the preposition “*lī*”. Note that in general if this invariant preceded in  $n - 1$  by the invariant “*fī*”, it will be followed in  $n + 1$  by “*'anna*”.

The Formal grammar that accepts phrases in which “*hīna*” plays the role of a temporal subordinating conjunction is presented in the form of a finite automaton comprising 12 states as the following figure (Fig.2) shows.  $G(hīna, CST) = \{N, \Sigma, P, q_0, q_{11}\}$ , with  $N = \{q_0, q_1, q_2, q_3, q_4, q_5, q_6, q_7, q_8, q_9, q_{10}, q_{11}\}$  is the not terminal vocabulary,  $\Sigma = \{VAC, VIN, GN, GV, CC, PD, PhN, hīna, 'inna, 'anna, 'an, 'ammā, fī, fa, \epsilon\}$  is the input vocabulary,  $P$  is the set of production rules and  $q_0 \in N$  is an initial state,  $q_{11} \in N$  is the final state.

- $G(hīna, CST) =$
1.  $q_0 \rightarrow CC\ q_1\ | \ PhN\ q_3\ | \ fī\ q_3\ | \ 'ammā\ q_3\ | \ \epsilon\ q_3,$
  2.  $q_1 \rightarrow VAC\ q_2\ | \ PD\ q_3\ | \ \epsilon\ q_3,$
  3.  $q_2 \rightarrow GN\ q_3,$
  4.  $q_3 \rightarrow hīna\ q_4,$
  5.  $q_4 \rightarrow VAC\ q_5\ | \ VIN\ q_5\ | \ 'anna\ q_8\ | \ \epsilon\ q_5,$
  6.  $q_5 \rightarrow GN\ q_6,$
  7.  $q_6 \rightarrow VAC\ q_8\ | \ VIN\ q_9\ | \ fa\ q_7\ | \ \epsilon\ q_{11},$
  8.  $q_7 \rightarrow 'inna\ q_8,$
  9.  $q_8 \rightarrow GN\ q_{11},$
  10.  $q_9 \rightarrow 'an\ q_{10},$
  11.  $q_{10} \rightarrow GV\ q_{11}.$

**Noun (hīn).** In general, the lexical invariant “*hīn*” occupies the function of a name in the case where it intervenes either in the middle or at the end of the sentence as shown in the following examples (3) and (4):

ويراك في كل حين وكل لحظة « *wa yarāka fī kulli hīn wa kulli laḥẓat.* » (He always saw you at all times). (3)

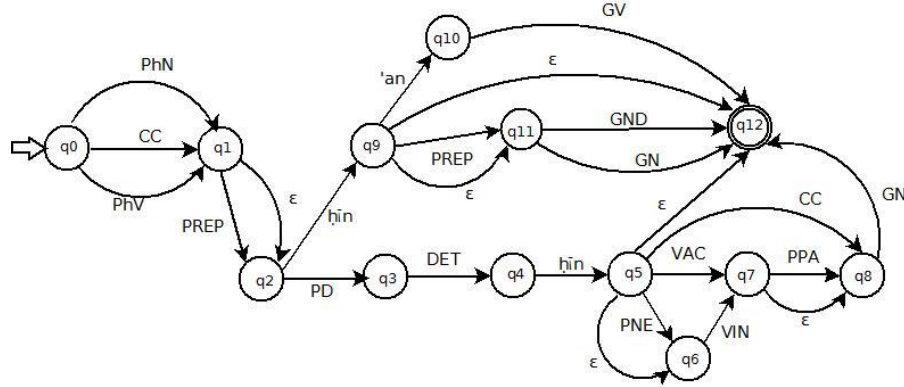


Fig. 3. Grammar diagram of "hīn" noun.

ولتعلن نبأه بعد حين « *wa lita 'lamanna naba'ahu ba 'da hīn.* » (You will soon hear about it). (4)

According to the two sentences above (3 and 4), we note that "hīn" expresses the notion of a short or long moment. That is why we can be replaced it by the word "waqt" (moment). So in this case, "hīn" has a syntactic and semantic equivalence with the word "waqt" as shown in the examples (3') and (4').

ويراك في كل وقت وكل لحظة « *wa yarāka fī kulli waqt wa kulli laḥẓat.* » (He always saw you at all times). (3')

ولتعلن نبأه بعد وقت « *wa lita 'lamanna naba'ahu ba 'da waqt.* » (You will soon hear about it). (4')

Note that in this case "hīn" can concatenate with the article "Al" to transform it from a noun to a determined noun as shown in the following sentence (5) or with the preposition "li" (example 6).

فحتى ذلك الحين أترككم في رعاية الله « *faḥattā ḍalika Alhīn 'atrukukum fī ri 'āyat allah.* » (And in the meantime, may God protect you). (5)

ينتظرون لحين ان يأتي دورهم في الحكاية « *yantaṭirūn liḥīn 'an ya 'tī dawruhum fī Alḥikāyat.* » (They will wait for a moment to come to their role in history). (6)

Similarly, in this sentence (5) if the invariant "hīn" is replaced by the word "waqt", the structure and meaning of the sentence do not move (Example 5'). Consequently, we can permute the sequence "liḥīn" with the lexical invariant "ilā" or "ḥattā" without any change in the source sentence (example 6').

فحتى ذلك الوقت أترككم في رعاية الله « *faḥattā ḍalika Alwaqt 'atrukukum fī ri 'āyat allah.* » (And until that time God bless you). (5')

ينتظرون إلى ان يأتي دورهم في الحكاية « *yantaṭirūn 'ilā 'an ya 'tī dawruhum fī Alḥikāyat.* » (They will wait for the time that will come their role in the history). (6')

In this case, the syntactic expectations of "hīn" are less complex than those of a subordinating conjunction. Indeed, it expects in  $n + 1$  either the empty set or a nominal group. This nominal group can be in the form "min GND" (example 7).

هل اتى على الانسان حين من الدهر « *hal 'atā 'alā Al'insān hīn min Aldahr.* » (Did the man go through a time lapse). (7)

The Formal grammar that accepts phrases in which "hīna" plays the role of a temporal subordinating conjunction is presented in the form of a finite automaton

**Table 6.** Linguistic Properties of "*hīn*".

<b>hīn / hīna</b>		
Graphic form / grammatical class	hīna / CST	hīn / Noun
Meaning	Temporal: <i>fī alwaqt allaḡī</i>	Expresses the notion of a moment
Removal	No	No
Position in the sentence	Beginning or middle	Middle or end
Change of position	Yes	No
Agglutination in n-1	<i>wa, fa</i>	<i>li, Al</i>
Agglutination in n+1	<i>mā, ḡāk, 'iḡan</i>	No
Form of sentence	<i>hīn P Q</i> ou <i>Q hīn P</i>	<i>P hīn</i> ou <i>P hīn Q</i>
Switching with other invariants	<i>'indamā</i> <i>fī hīna 'anna = baynamā</i> <i>lihīn = ḡattā</i>	<i>waqt</i> <sup>55</sup> <i>lihīn = 'ilā</i>
Syntactic expectations at n + 1	GV or GN	GN, min GND or Ø
Governance	Governs two kernels: verb + subject and / or theme + predicate.	No
Discontinuous	No	No
Traduction	When	Moment

comprising 13 states as the following figure (Fig.3) shows.  $G(hīn, \text{Noun}) = \{N, \Sigma, P, q_0, q_{12}\}$ , with,  $N = \{q_0, q_1, q_2, q_3, q_4, q_5, q_6, q_7, q_8, q_9, q_{10}, q_{11}, q_{12}\}$  is the not terminal vocabulary,  $\Sigma = \{VAC, VIN, GN, GND, GV, CC, PD, DET, PREP, PNE, PPA, PhN, PhV, hīn, 'an, \varepsilon\}$  is the input vocabulary,  $P$  is the set of production rules and  $q_0 \in N$  is an initial state,  $q_{12} \in N$  is the final state.

$$G(hīn, \text{Noun}) = \left\{ \begin{array}{l} 1) \quad q_0 \rightarrow PhN \ q_1 | PhV \ q_1 | CC \ q_1, \\ 2) \quad q_1 \rightarrow PREP \ q_2 | \varepsilon \ q_2, \\ 3) \quad q_2 \rightarrow PD \ q_3 | hīn \ q_9, \\ 4) \quad q_3 \rightarrow DET \ q_4, \\ 5) \quad q_4 \rightarrow hīn \ q_5, \\ 6) \quad q_5 \rightarrow VAC \ q_7 | PNE \ q_6 | CC \ q_8 | \varepsilon \ q_6 | \varepsilon \ q_{12}, \\ 7) \quad q_6 \rightarrow VIN \ q_7, \\ 8) \quad q_7 \rightarrow PPA \ q_8 | \varepsilon \ q_8, \\ 9) \quad q_8 \rightarrow GN \ q_{12}, \\ 10) \quad q_9 \rightarrow PREP \ q_{11} | 'an \ q_{10} | \varepsilon \ q_{11} | \varepsilon \ q_{12}, \\ 11) \quad q_{11} \rightarrow GV \ q_{12}. \end{array} \right.$$

The different linguistic properties of "*hīn*/ *hīna*" are summarized in the table below:

<sup>55</sup> In this case, the word "waqt" does not designate a lexical invariant but rather a noun.

**Table 7.** Linguistic Properties of " 'aṭnā'a ".

'aṭnā'a	
Graphic ambiguity	Unique : إنشاء « 'aṭnā'a », إنشاء « aṭnā'a »
Grammatical ambiguity	Unique: Preposition
Meaning	Temporal: ḥilāla
Removal	No
Position in the sentence	Beginning or middle
Change of position	Yes
Agglutination in n-1	wa, fa, ka
Agglutination in n+1	Attached personal pronoun: h, hā...
Form of sentence	'aṭnā'a P Q ou Q 'aṭnā'a P
Switching with other invariants	Yes: ḥilāla
Syntactic expectations at n + 1	GND
Governance	No
Discontinuous	No
Traduction	During

#### 4.2 The Lexical Invariant " 'aṭnā'a "

According to Arabic literature, the word " 'aṭnā'a " is a preposition that connects two or more events in the sense of "ḥilāla". These events are parallel and simultaneous with respect to time. If we take into account how to transcribe the "Hamza", the lexical invariant " 'aṭnā'a " is not graphically ambiguous.

The aim of our study is to have more or less precise idea on this lexical invariant. To do this, we first selected some samples of sentences that contain " 'aṭnā'a " from our corpus [5, 6]. Our study sample of " 'aṭnā'a " contains 50 sentences. We used this sample at first to identify the notable features. Among these pertinent remarks, we found that " 'aṭnā'a " was always followed by a determined nominal group.

Indeed, to better understand the functioning of this invariant in the Arabic language, we have treated some selected examples from our study sample.

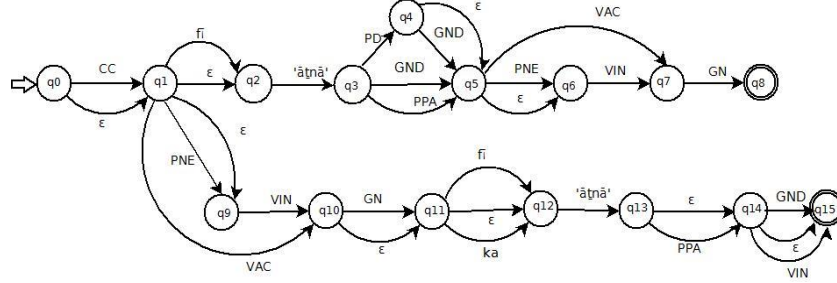
Consider the following examples:

« wa 'aṭnā'a dirāsati fī brīṭāniyā qaddamnā 'amalāa kuntu 'a 'mal fih musā'id muḥrig ». (During my studies in the United Kingdom, we presented a work in which I was an assistant director) (7).

« ya 'ḥuḍu alriwāyat ma 'ahu 'aṭnā'a safarihi illā alḥārig » (He took the novel with him during his travels abroad). (8)

From examples above, we note that " 'aṭnā'a " can intervene either at the beginning or in the middle of the sentence and never at the end of the sentence except in the case where it is agglutinated in n + 1 with a Personal pronoun attached as shown in the following sentence:

« faliyahuṣ almuslim alnabīh alnaziḥ alā muḥāsabat naḥsihi almuḥāsabat alšāmilat qabla al'amal wa 'aṭnā'ahu » (The true Muslim must be vigilant in asking for complete accounts before and during the work). (9)



**Fig. 4.** Grammar diagram of “‘atnā'a” subordinating conjunction.

As we noted above, this invariant takes the meaning of “*hīlāla*”. That is, we can switch “‘atnā'a” by “*hīlāla*” and keeping the same meaning and structure of the sentence as shown in the following example:

يأخذ الرواية معه أثناء سفره إلى الخارج ↔ يأخذ الرواية معه خلال سفره إلى الخارج (He took the novel with him during his travels abroad).

The table below summarizes these different linguistics properties:

The Formal grammar that accepts phrases in which “‘atnā'a” plays the role of a temporal subordinating conjunction is presented in the form of a finite automaton comprising 16 states as the following figure (Fig.4) shows.  $G('atnā'a, CS) = \{N, \Sigma, P, q_0, q_{15}\}$ , with  $N = \{q_0, q_1, q_2, q_3, q_4, q_5, q_6, q_7, q_8, q_9, q_{10}, q_{11}, q_{12}, q_{13}, q_{14}, q_{15}\}$  is the not terminal vocabulary,  $\Sigma = \{VAC, VIN, GN, CC, PPA, PNE, PD, fi, ka, 'atnā'a, \epsilon\}$  is the input vocabulary,  $P$  is the set of production rules and  $q_0 \in N$  is an initial state,  $q_{15} \in N$  is the final state.

$$G('atnā'a, CS) = \left\{ \begin{array}{l} 1) \quad q_0 \rightarrow CC \ q_1 \mid \epsilon \ q_1, \\ 2) \quad q_1 \rightarrow fi \ q_2 \mid VAC \ q_{10} \mid PNE \ q_9 \mid \epsilon \ q_2 \mid \epsilon \ q_9, \\ 3) \quad q_2 \rightarrow 'atnā' \ q_3, \\ 4) \quad q_3 \rightarrow PD \ q_4 \mid GND \ q_5 \mid PPA \ q_5, \\ 5) \quad q_4 \rightarrow GND \ q_5 \mid \epsilon \ q_5, \\ 6) \quad q_5 \rightarrow PNE \ q_6 \mid VAC \ q_7 \mid \epsilon \ q_6, \\ 7) \quad q_6 \rightarrow VIN \ q_7, \\ 8) \quad q_7 \rightarrow GN \ q_8, \\ 9) \quad q_9 \rightarrow VIN \ q_{10}, \\ 10) \quad q_{10} \rightarrow GN \ q_{11} \mid \epsilon \ q_{11}, \\ 11) \quad q_{11} \rightarrow fi \ q_{12} \mid ka \ q_{12} \mid \epsilon \ q_{12}, \\ 12) \quad q_{12} \rightarrow 'atnā' \ q_{13}, \\ 13) \quad q_{13} \rightarrow PPA \ q_{14} \mid \epsilon \ q_{14}, \\ 14) \quad q_{14} \rightarrow GND \ q_{15} \mid VIN \ q_{15} \mid \epsilon \ q_{15}. \end{array} \right.$$

## 5 Conclusions and Perspectives

In this work, we tried to study and classify 13 simple Arabic lexical invariants. In our study, we have focused on the notions of “permutation” or “commutation” that based the approach of structural linguistics called “distributional”. To do this, first for each lexical invariant, we developed a sample corpus that contains different contexts of the

lexical invariant in question. Second, from this corpus, we identified the different linguistic criteria of the lexical invariant that allow us to correctly identify it.

Finally, we codified this information in the form of linguistic rules in order to model it by diagram of grammar (finite state machine). Indeed, as perspectives for this work, we will evaluate our temporal invariant grammars on a big data.

## **6 Appendix**

VIN: Verb present tense, VAC: verb past tense, VI: imperative verb, N: Noun, PREP: preposition, CC: coordinating conjunction, CS: subordinating conjunction, ADVT: adverb of time, PNE: negation particle, CCOND: conditional conjunction, CST: temporal subordinating conjunction, PS: Particle of surprise, PD: demonstrative pronoun, PF: Future particle, PPA: Pronoun personal attached. GN: Nominal group, GND: Determined nominal group, GV: verbal group, PhN: Nominal sentence, PhV: verbal sentence, DET: determiner.

## **References**

1. Jaccarini, A., Gaubert, C.: Le programme Mogador en linguistique formelle arabe et ses applications dans le domaine de la recherche et du filtrage sémantique (2012)
2. Audebert, C., Jaccarini, A.: À la recherche du khabar, outils en vue de l'établissement d'un programme d'enseignement assisté par ordinateur. Institut français d'archéologie orientale du Caire, pp. 217–256 (1986)
3. Ghoul, D.: Classifications et grammaires des invariants lexicaux arabes en prévision d'un traitement informatique de cette langue les invariants temporels. Doctorat (2016)
4. Jaccarini, A.: Grammaires modulaires de l'arabe. Mise en œuvre informatique et stratégies. Thèse de doctorat, Sorbonne (1997)
5. Ghoul, D.: Construction d'un corpus arabe à partir du Web dans le but d'identifier les mots-outils ou tokens. JADT: Journées internationales d'Analyse statistique des Données Textuelles, INALCO, pp. 271–276 (2014)
6. Ghoul, D., Ibrahim, A. H.: Web arabic corpus : construction d'un large corpus arabe annoté morpho syntaxiquement à partir du Web. CECTAL'15, pp. 12–16 (2015)
7. Gaubert, C.: Kawâkib, une application pour le traitement automatique de textes arabes. no. 44, pp. 66–81 (2010)
8. Ghoul, D., Ibrahim, A. H., Audebert, C.: Rules-based grammatical and semantic disambiguation of the token "hatta". In: 5th International Conference on Information Communication Technology and Accessibility (ICTA), pp. 1–6 (2015)
9. Dubois, J., Dubois-Charlier, F.: Principes et méthode de l'analyse distributionnelle, Langages, vol. 5, no. 20, pp. 3–13 (1970)
10. Dubois, J.: Grammaire distributionnelle, vol. 1, no. 1, pp. 41–48 (1969)